

NOTE

USING *DAUBERT* TO EVALUATE EVIDENCE-BASED SENTENCING

Charlotte Hopkinson[†]

INTRODUCTION	724
I. BACKGROUND TO EVIDENCE-BASED SENTENCING	725
II. SCIENTIFIC VALIDITY OF PREDICTIVE INSTRUMENTS DURING SENTENCING	729
A. Capital Cases	730
B. Sexually Violent Persons	732
C. Analogy to Evidence-Based Sentencing	733
III. <i>DAUBERT</i> 'S APPLICATION TO EVIDENCE-BASED SENTENCING GENERALLY AND IN PENNSYLVANIA	733
A. Introduction to the Pennsylvania Commission on Sentencing Risk Assessment Tool	736
B. Prong 1: Whether Evidence-Based Sentencing Has Been Tested?	737
C. Prong 2: Whether Evidence-Based Sentencing Has Been Subjected to Peer Review and Publication?	742
D. Prong 3: What Is the Known (or Potential) Error Rate and Are There Standards that Control Evidence-Based Sentencing?	744
E. Prong 4: Whether Evidence-Based Sentencing Is Generally Accepted Within a Relevant Scientific Community and to What Degree? ...	746
IV. <i>DAUBERT</i> 'S INTERACTION WITH FEDERAL RULE OF EVIDENCE 403	748
A. Sex	749
B. Age	751

[†] Charlotte Hopkinson is a student in the dual degree Cornell Law School and Université Paris I Panthéon-Sorbonne program and will graduate in 2019 with a JD and Master en Droit. Special thanks to all those who were willing to discuss my Note with me, Joe Margulies and John Zipp. John Zipp was invaluable in interpreting the methods and results published by the Pennsylvania Commission on Sentencing. This Note was inspired by my Science and the Law Class at Auburn Correctional Facility, which was taught through the Cornell Prison Education Program in Fall 2016. Student-led class discussions helped me consider various legal and policy viewpoints.

V. PUBLIC POLICY CONCERNS	753
CONCLUSION	754

INTRODUCTION

Jack and Jill went up the hill,
to steal a pail of water,
Both were caught and sentenced to jail,
But Jack came out two years later.

Why? Assume that both Jack and Jill's cases are identical in facts, procedure, jury composition, and verdict. The only relevant difference is that Jack is a man and Jill is a woman. Statistically, men are more likely to recidivate than women,¹ and under a sentencing system called evidence-based sentencing, the judge agrees that Jack is more likely to commit another crime due to his gender. The solution to this prediction? Jack will spend more time in prison, simply because he is a man.

Evidence-based sentencing is part of a growing trend of using actuarial risk assessments in the criminal justice system.² Evidence-based sentencing relies on a set of factors, some of which are relevant to the crime committed but most of which are based on immutable characteristics outside of the defendant's control, to predict the probability that a defendant will recidivate.³ Notwithstanding the constitutionally problematic nature of relying on certain factors, such as sex,⁴ the scientific validity of such methods is questionable. Although evidentiary rules do not apply to sentencing, this Note argues that the admissibility of such evidence in general, and specifi-

¹ Sonja B. Starr, *Evidence-Based Sentencing and the Scientific Rationalization of Discrimination*, 66 STAN. L. REV. 803, 863 (2014) ("To give a simple illustration, if a sentence is based on crime severity plus gender, and these factors together produce a ten-year sentence for a male when an otherwise-identical woman would have received seven years, male gender is not solely responsible for the sentence; crime severity establishes the baseline of seven years. But male gender is solely responsible for the extra three years.").

² Dawinder S. Sidhu, *Moneyball Sentencing*, 56 B.C. L. REV. 671, 673 (2015) ("More and more, courts today are adopting the use of risk-assessment tools in sentencing. These risk-assessment tools take information on recidivism rates for groups and use them to estimate the risk of recidivism for individuals possessing those same group characteristics.").

³ The topic of predictive factors of dangerousness recently arose in oral arguments in front of the Supreme Court in *Buck v. Davis*. Transcript of Oral Argument at 31, *Buck v. Davis*, 137 S. Ct. 759 (2016) (No. 15-8049).

⁴ The author of this Note would like to point out that actuarial instruments and researchers tend to use gender and sex interchangeably. These two words, however, have two different meanings. For the purpose of this Note, sex is biologically defined, whereas gender is a societal construct. The author also uses the singular "they" throughout this Note.

cally Pennsylvania's new sentencing scheme, fails the *Daubert* framework.

Part I of this Note provides a background into the characteristics and history of different types of evidence-based sentencing and addresses the arguments for and against its use. Part II describes the "future dangerousness" standard used in death penalty sentencing and civil commitments of sexually violent persons as an analogy to evidence-based sentencing and argues that evidence-based sentencing must withstand some type of evidentiary gatekeeping test. Part III subjects evidence-based sentencing to the *Daubert* test and applies this test broadly to evidence-based sentencing generally, as well as specifically to Pennsylvania's Risk Assessment Project.⁵ Part IV uses Federal Rule of Evidence 403 to further explain why evidence-based sentencing should not be used, particularly for models that include sex and age. Part V briefly addresses certain penological theories that evidence-based sentencing reflects and questions whether evidence-based sentencing is the appropriate solution.

I

BACKGROUND TO EVIDENCE-BASED SENTENCING

Evidence-based sentencing focuses on predicting an individual's recidivism risk based on empirical research.⁶ Evidence-based sentencing is a type of risk assessment, or actuarial method, that relies on a large dataset to evaluate the "statistical correlations between a group trait and that group's criminal offending rate" as opposed to a clinical evaluation.⁷ Empirical research attributes recidivism to a wide range of factors, including criminal history, sex, age, marital status, employment, education, parental convictions, family members who have been crime victims, high school grades, chances of finding work above minimum wage, dependence on social as-

⁵ See PA. COMM'N ON SENT'G, OVERVIEW OF THE RISK ASSESSMENT INSTRUMENT, <http://pcs.la.psu.edu/publications-and-research/research-and-evaluation-reports/risk-assessment> [<https://perma.cc/5FUC-GVX2>].

⁶ Starr, *supra* note 1, at 805.

⁷ Bernard H. Harcourt, *Against Prediction: Sentencing, Policing, and Punishing in an Actuarial Age* 11 (U. of Chi. Pub. Law & Legal Theory Working Paper Series, Paper No. 94, 2005); see Jordan M. Hyatt, Mark H. Bergstrom & Steven L. Chanenson, *Follow the Evidence: Integrate Risk Assessment into Sentencing*, 23 FED. SENT'G REP. 266, 266 (2011) [hereinafter *Follow the Evidence*] (noting that actuarial risk assessments rely on static variables whereas clinical risk assessments rely on dynamic variables).

sistance, finances, and neighborhood of residence.⁸ By using an impersonal, data-driven method to determine the recidivism risk, many actors in the criminal justice system believe that they are simultaneously increasing fairness in discretionary sentencing and reducing the overpopulation of prisons by diverting low-risk offenders from prison.⁹

Although probabilistic reasoning is found throughout the criminal justice system, such as a police officer's determination that there is sufficient cause for a search, forensic testimony concerning a DNA "match," or the legal standard for criminal guilt itself ("beyond a reasonable doubt"), these examples of probabilistic reasoning are not based on statistical correlations of a group trait, but rather on a situational case-by-case analysis.¹⁰ This new focus on evidence-based sentencing changes the relationship between group data and individual predictions in the criminal justice system and results in serious consequences (loss of liberty). An increasing amount of states are adopting evidence-based sentencing guidelines,¹¹ and the American Law Institute's Model Penal Code¹² draft incorpo-

⁸ These factors differ between states and assessment methods. Researchers have identified hundreds of factors relevant to sentencing but are trying to narrow them down to the few factors that are the most predictive. See Starr, *supra* note 1, at 805, 811–13; Sonja B. Starr, *Sentencing, by the Numbers*, N.Y. TIMES (Aug. 10, 2014), <http://www.nytimes.com/2014/08/11/opinion/sentencing-by-the-numbers.html> [<https://perma.cc/EML7-TRAY>]; see also J.C. Oleson, *Risk in Sentencing: Constitutionally Suspect Variables and Evidence-Based Sentencing*, 64 S.M.U. L. REV. 1329, 1350–51 (2011) (distinguishing "having criminal companions," "antisocial personality," "criminogenic needs," "adult criminal history," and "race" as strongly significant predictors of recidivism versus "substance abuse," "family structure," "intellectual functioning," "family criminality," "gender," "socio-economic status of origin," and "personal distress" as weaker, but still significant, predictors). For a list of 125 factors, see PA. COMM'N ON SENT'G, INTERIM REPORT 1: REVIEW OF FACTORS USED IN RISK ASSESSMENT INSTRUMENTS 13–17 (2011).

⁹ Starr, *supra* note 8; see Julia Angwin, Jeff Larson, Surya Mattu & Lauren Kirchner, *Machine Bias*, PROPUBLICA (May 23, 2016), <https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing> [<https://perma.cc/8ZPQ-9F92>]; *Follow the Evidence*, *supra* note 7, at 267 (arguing that the sentences would be facially disparate, but individualized and thus fair). The use of evidence-based models is also used in bail determinations. Shaila Dewan, *Judges Replacing Conjecture with Formula for Bail*, N.Y. TIMES (June 26, 2015), <http://www.nytimes.com/2015/06/27/us/turning-the-granting-of-bail-into-a-science.html> [<https://perma.cc/5X6J-5KD8>].

¹⁰ Harcourt, *supra* note 7, at 11.

¹¹ See NANCY LAVIGNE ET AL., URBAN INST., JUSTICE REINVESTMENT INITIATIVE STATE ASSESSMENT REPORT 39–40 (2014), <https://www.urban.org/sites/default/files/publication/22211/412994-Justice-Reinvestment-Initiative-State-Assessment-Report.PDF> [<https://perma.cc/74XL-3N3K>]; see also Starr, *supra* note 8 (noting that "this practice has rapidly expanded much more recently").

¹² MODEL PENAL CODE: SENTENCING app. A at 133, 135 (AM. LAW INST., Discussion Draft No. 4, 2012).

rates a recommendation to use actuarial instruments.¹³ The use of predictive instruments to estimate recidivism originated with discretionary parole (which, ironically, has been abolished in sixteen states)¹⁴ but has expanded to more than twenty states' current sentencing practices.¹⁵ While Pennsylvania is debuting its own predictive instrument,¹⁶ other states use generic actuarial tests that are only sometimes calibrated to state-specific populations. Two of the most common generic actuarial tests are the Level of Services Inventory-Revised (LSI-R) and the Correctional Offender Management Profiling for Alternative Sanctions (COMPAS).¹⁷ Like many of the state-specific tests, these commercial tests include objective and subjective variables that are unrelated to the defendant's crime.¹⁸

Both legal and scholarly advocates for evidence-based sentencing believe in its ability to simultaneously hold offenders accountable, protect the public, and reduce incarceration (and associated financial costs and social harms).¹⁹ Evidence-based sentencing purportedly offers a "scientific" way to systematically differentiate low-risk and high-risk offenders. Advocates argue that judges already engage in predictions during sentencing—therefore, evidence-based sentencing will guide judges with scientific feedback and increase sentencing transparency.²⁰ The purported benefits of this transparency extend

¹³ See Starr, *supra* note 1, at 805.

¹⁴ TIMOTHY A. HUGHES ET AL., U.S. DEPT OF JUSTICE, BUREAU OF JUSTICE STATISTICS, TRENDS IN STATE PAROLE, 1990–2000 (2001). In 1994, Virginia abolished parole and established a risk assessment instrument for reducing incarceration rates. David A. Soule & Stacy S. Najaka, *Paving a Path to Informed Sentencing Decisions*, 25 FED. SENT'G REP. 181, 181 (2013).

¹⁵ Starr, *supra* note 1, at 809 ("A review of case law, sentencing commission websites, and relevant legislation indicates that at least twenty states' courts are now, in some or all cases, incorporating actuarial risk assessments into the determination of the defendant's sentence.").

¹⁶ Logan Koepke, *Pennsylvania Will Vary Jail Terms for the Same Crime, Based on Where You Live*, EQUAL FUTURE (Sept. 16, 2015), <https://www.equalfuture.us/2015/09/16/pennsylvania-will-vary-jail-terms-for-the-same-crime-based-on-where-you-live/> [<https://perma.cc/Z7HU-UQMN>].

¹⁷ Starr, *supra* note 1, at 812.

¹⁸ *Id.* at 812–13.

¹⁹ See *id.* at 816; Soule & Najaka, *supra* note 14, at 181; see also *Follow the Evidence*, *supra* note 7, at 266 ("The use of risk assessment at sentencing underscores an overall shift in the purposes of sentencing, moving from a backward-looking retributive approach with a focus on uniformity, proportionality, and reduction of unwarranted disparity to an approach that also incorporates a formalized, forward-looking utilitarian goal.").

²⁰ See Lynn S. Branham, *Follow the Leader: The Advisability and Propriety of Considering Cost and Recidivism Data at Sentencing*, 24 FED. SENT'G REP. 169, 171 (2012).

to the general public who will perceive the legal system as more legitimate, thus causing them to comply with more laws.²¹

Evidence-based sentencing has also come under scrutiny and criticism.²² Scholars note that because evidence-based sentencing focuses on uncontrollable characteristics of the defendant rather than the crime at hand, it distorts penological goals, particularly the retributive and rehabilitative goals.²³ Other critics, particularly in the media, analogize evidence-based sentencing to the movie *Minority Report* in which the government punishes crimes before they happen, because a person's prison sentence is partially based on their future likelihood to recidivate.²⁴ Furthermore, the constitutionality of the evidence-based sentencing is at question.²⁵ The use of demographic factors including sex and socioeconomic status raises constitutional issues such as equal protection and the rights of indigent defendants.²⁶ Although these scholars support reducing incarceration rates, they predict that doing so by considering demographic factors risks further exacerbating the inequities of the criminal justice system.²⁷ These critiques are not limited to academics. Although some courts have apprehensively embraced evidence-based sentencing,²⁸ other actors

²¹ See Tracey L. Meares, *Three Objections to the Use of Empiricism in Criminal Law and Procedure—And Three Answers*, 2002 U. ILL. L. REV. 851, 866 (2002).

²² See Starr, *supra* note 1, at 817 (“Although most of the literature on EBS is positive, or even celebratory, a few scholars have criticized it.”).

²³ See Harcourt, *supra* note 7, at 31–32 (“We have come to associate the prediction of future criminality with just punishment.”); see also Starr, *supra* note 1, at 817 (“[P]rediction instruments contravene punishment theory, because punishment turns on *who the defendant is . . .*”).

²⁴ Jared Greenhouse, *Pennsylvania Wants to Use Science in Criminal Sentencing*, HUFFINGTON POST (Aug. 11, 2015, 11:54 AM), http://www.huffingtonpost.com/entry/science-criminal-sentencing_us_55c8ed49e4b0923c12bda693 [https://perma.cc/U272-PUS8]; see Anna Maria Barry-Jester, Ben Casselman & Dana Goldstein, *The New Science of Sentencing*, THE MARSHALL PROJECT (Aug. 4, 2015, 7:15 AM), [hereinafter *The New Science of Sentencing*] <https://www.themarshallproject.org/2015/08/04/the-new-science-of-sentencing#.4mF2KXDVR> [https://perma.cc/9465-4MND].

²⁵ See Starr, *supra* note 1, at 821–42; Oleson, *supra* note 8, at 1372–88; Shaina D. Massie, Note, *Orange Is the New Equal Protection Violation: How Evidence-Based Sentencing Harms Male Offenders*, 24 WM. & MARY BILL RTS. J. 521, 522 (2015) (“Penological considerations of gender in sentencing are simply incompatible with abstract notions that criminal offenders appear before the court in their individual capacities. More important, the use of gender in evidence-based sentencing violates the concrete promises of equal protection under the law provided by the Constitution.”).

²⁶ See Starr, *supra* note 1, at 806.

²⁷ See *id.*

²⁸ Danielle Citron, *(Un)Fairness of Risk Scores in Criminal Sentencing*, FORBES (July 13, 2016, 3:26 PM), <http://www.forbes.com/sites/daniellecitron/2016/07/13/unfairness-of-risk-scores-in-criminal-sentencing/#46d161f54479> [https://

in the justice system, such as former Attorney General Eric Holder, have warned of the dangers of using immutable characteristics in sentencing decisions.²⁹

Both sides of the evidence-based sentencing debate have laudable goals; however, it is difficult, ironically, to predict the outcome of evidence-based sentencing in the criminal justice system. Rather than focus on the hypothetical positive and negative outcomes or the constitutionality of this sentencing practice, this Note seeks to further contribute to academic discussion by exploring the scientific validity of such instruments.

II

SCIENTIFIC VALIDITY OF PREDICTIVE INSTRUMENTS DURING SENTENCING

Much of the legal and academic discussion surrounding the validity and use of predictive instruments is focused on two areas of the law: future dangerousness predictions for capital

perma.cc/H4DJ-FQYB] (“[The Wisconsin Supreme Court] made clear its concerns about the accuracy and potential bias of risk scoring systems. The Court recognized that while states like New York have studied the effectiveness and predictive accuracy of COMPAS scores and found them ‘satisfactory,’ Wisconsin had not completed a statistical validation study for COMPAS for its population.”).

²⁹ Attorney General Holder emphasized that he was “concerned that [risk assessments] may inadvertently undermine our efforts to ensure individualized and equal justice. By basing sentencing decisions on static factors and immutable characteristics—like the defendant’s education level, socioeconomic background, or neighborhood—they may exacerbate unwarranted and unjust disparities that are already far too common in our criminal justice system and in our society.” Eric Holder, Attorney General, Speech at the National Association of Criminal Defense Lawyers 57th Annual Meeting and 13th State Criminal Justice Network Conference (Aug. 1, 2014) (transcript available at <https://www.justice.gov/opa/speech/attorney-general-eric-holder-speaks-national-association-criminal-defense-lawyers-57th> [<https://perma.cc/8J5G-FFKU>]); see also Letter from Jonathan J. Wroblewski, Director, Office of Policy and Legislation, U.S. Department of Justice, to Hon. Patti B. Saris, Chair, U.S. Sentencing Commission, at 7 (July 29, 2014), <https://www.justice.gov/sites/default/files/criminal/legacy/2014/08/01/2014annual-letter-final-072814.pdf> [<https://perma.cc/WMP8-ZGUV>] (“[D]etermin[ing] risk levels based on static, historical offender characteristics such as education level, employment history, family circumstances and demographic information . . . [and] basing criminal sentences, and particularly imprisonment terms, primarily on such data—rather than the crime committed and surrounding circumstances—is a dangerous concept that will become much more concerning over time as other far reaching sociological and personal information unrelated to the crimes at issue are incorporated into risk tools. This phenomenon ultimately raises constitutional questions because of the use of group-based characteristics and suspect classifications in the analytics. Criminal accountability should be primarily about prior bad acts proven by the government before a court of law and not some future bad behavior predicted to occur by a risk assessment instrument.”).

cases³⁰ and for sexually violent person civil commitments.³¹ This section reviews the relevant legal and academic discourse and analogizes the use of future dangerousness predictions in these areas to the use of evidence-based predictions at sentencing. Although courts have found the use of future dangerousness predictions in capital cases and civil commitments admissible, many academics argue that a gatekeeping test like *Daubert* would disallow this kind of evidence.³²

A. Capital Cases

In the sentencing stages of capital cases, experts may testify on a capital defendant's future dangerousness.³³ In *Barefoot v. Estelle*,³⁴ the United States Supreme Court refused to categorically exclude a psychiatrist's testimony about the defendant's future dangerousness in sentencing, despite the American Psychiatric Association's assertion that these kinds of predictions are wrong *most* of the time.³⁵ In this case, the Supreme Court echoed the views of some proponents of evidence-based sentencing—notably that some level of prediction is already present in all aspects of the criminal justice system, such as bail decisions and parole hearings.³⁶ A similar case arose in *United States v. Fields*,³⁷ where the appellant did not argue that “psychiatric predictions of future dangerousness during the punishment phase are inadmissible *per se*” but challenged whether the expert opinion offered was reliably sufficient to be introduced at sentencing.³⁸ In the time between deciding *Barefoot* and *Fields*, the United States Supreme Court had ruled on *Daubert v. Merrell Dow Pharmaceuticals, Inc.*,³⁹ in which it interpreted Federal Rule of Evidence 702⁴⁰ and guided

³⁰ See Brian Sites, Note, *The Danger of Future Dangerousness in Death Penalty Use*, 34 FLA. ST. U. L. REV. 959, 967–69 (2007).

³¹ See Eric S. Janus & Robert A. Prentky, *Forensic Use of Actuarial Risk Assessment with Sex Offenders: Accuracy, Admissibility and Accountability*, 40 AM. CRIM. L. REV. 1443, 1443–44 (2003).

³² See *infra* subpart II.A (Capital Cases) and subpart II.B (Sexually Violent Persons).

³³ See Erica Beecher-Monas & Edgar Garcia-Rill, *Danger at the Edge of Chaos: Predicting Violent Behavior in a Post-Daubert World*, 24 CARDOZO L. REV. 1845, 1895–1901 (2003).

³⁴ 463 U.S. 880 (1983).

³⁵ See *id.* at 899–904.

³⁶ See *id.* at 897.

³⁷ 483 F.3d 313 (5th Cir. 2007).

³⁸ *Id.* at 341.

³⁹ 509 U.S. 579 (1993). *Daubert* supplemented Federal Rule of Evidence 702, which had superseded *Frye v. United States*, 293 F. 1013 (D.C. Cir. 1923).

⁴⁰ FED. R. EVID. 702.

courts acting as “gatekeepers” of expert testimony.⁴¹ In *Fields*, the Fifth Circuit ruled that the Federal Death Penalty Act did not require a *Daubert* analysis, despite the American Psychological Association’s urging that some sort of *Daubert*-like investigation should be carried out with regard to expert testimony in sentencing hearings.⁴²

Even in the few states that have called for some gatekeeping test for expert future dangerousness predictions in death penalty cases, most have refused to categorically exclude future dangerousness predictions and many “exempt predictive expertise from any special testing, applying instead other rules on expert testimony and standard tests for relevance.”⁴³ Only one judge has rejected predictive testimony under *Daubert* in a specially concurring opinion.⁴⁴ Most scholars agree that gatekeeping with regards to expert future dangerousness predictions in death penalty cases is necessary,⁴⁵ but disagree as to which standard to apply (that is, *Daubert* or a more relaxed standard) and whether future dangerousness predictions do, in

⁴¹ See *Daubert*, 509 U.S. at 600–601 (noting that even though “Rule 702 confides to the judge some gatekeeping responsibility in deciding questions of the admissibility of proffered expert testimony[,]” it does not “impose[] on them either the obligation or the authority to become amateur scientists in order to perform that role”).

⁴² *Fields*, 438 F.3d at 342.

⁴³ Alexander Scherr, *Daubert & Danger: The “Fit” of Expert Predictions in Civil Commitments*, 55 HASTINGS L.J. 1, 61–62, 73–75 (2003) (footnote omitted); see Sites, *supra* note 30, at 970–71 (noting that when courts are “faced with choosing between abandoning a desirable sentencing tool (future dangerousness) and relying on unsteady analytical tools (psychological testing), courts prefer the latter”).

⁴⁴ See *Flores v. Johnson*, 210 F.3d 456, 458–70 (5th Cir. 2000) (Garza, J., specially concurring) (subjecting the predictive testimony to a *Daubert* analysis and arguing that *Barefoot* should be overturned). But see *People v. Murtishaw*, 631 P.2d 446, 470 (Cal. 1981) (rejecting predictive dangerousness testimony not under *Kelly/Frye* but rather because the court found that the testimony’s prejudicial impact outweighed its probative value). The court in *Murtishaw* noted that in places where the jury is required by law to determine whether someone is “dangerous,” such predictive testimony, although unreliable, would often be the only evidence available to assist the trier of fact. *Id.* at 469; see 18 U.S.C. § 3593(c) (2012) (“Information is admissible regardless of its admissibility under the rules governing admission of evidence at criminal trials except that information may be excluded if its probative value is outweighed by the danger of creating unfair prejudice, confusing the issues, or misleading the jury.”).

⁴⁵ See generally Sites, *supra* note 30, at 987, 992–95 (arguing that *Daubert* should apply to death penalty sentencing proceedings and supplementary methods to minimize jury confusion); Erica Beecher-Monas, *The Epistemology of Prediction: Future Dangerousness Testimony and Intellectual Due Process*, 60 WASH. & LEE L. REV. 353, 364–79 (2003) (arguing that *Barefoot* was wrongly decided and proposing the constitutionalization of the *Daubert* test for death penalty sentencing proceedings).

fact, satisfy *Daubert*.⁴⁶ Despite the fact that *Daubert* does not apply to state or federal⁴⁷ sentencing procedures, academic and legal discussions have made it clear that *Daubert* offers a helpful framework through which to analyze the reliability of future dangerousness predictions in capital penalty sentencing.⁴⁸

B. Sexually Violent Persons

Predictive assessments of future dangerousness have also been used in civil commitment hearings of sexually violent persons. Civil commitments are used to continue to incapacitate sex offenders who have finished their criminal sentences.⁴⁹ Risk assessment plays a central role in this assessment⁵⁰ and is often mandated by state law.⁵¹

Despite the known flaws of such predictive data techniques, courts have continued to allow this kind of evidence.⁵² The legal and scholarly debate surrounding civil commitments highlights the lack of consensus over what standard of evidentiary admissibility should be used for expert testimony in civil commitment.⁵³ As in capital cases, some scholars have used *Daubert* as a framework to analyze the admissibility of future dangerousness predictions in sexual violent person civil commitment hearings.⁵⁴

⁴⁶ See generally Beecher-Monas & Garcia-Rill, *supra* note 33, at 1847, 1856–57, 1880, 1900 (applying *Daubert* as a constitutional ground to decide whether expert testimony on a specific defendant’s future dangerousness violates due process and concluding that clinical predictions must be excluded and that actuarial instruments “barely squeak through” *Daubert*).

⁴⁷ 18 U.S.C. § 3661 (2012); U.S. SENTENCING GUIDELINES MANUAL § 6A1.3, Commentary (U.S. SENTENCING COMM’N 2016) (“In determining the relevant facts, sentencing judges are not restricted to information that would be admissible at trial.”).

⁴⁸ See Beecher-Monas & Garcia-Rill, *supra* note 33, at 1859.

⁴⁹ See Janus & Prentky, *supra* note 31, at 1443.

⁵⁰ See *id.* at 1443–44.

⁵¹ See George G. Woodworth & Joseph B. Kadane, *Expert Testimony Supporting Post-Sentence Civil Incarceration of Violent Sexual Offenders*, 3 LAW, PROBABILITY AND RISK 221, 236 (2004). For example, “Virginia’s Sexually Violent Predator statute not only mandates the use of a specific instrument but also specifies the cutoff score on that instrument that must be achieved to proceed further in the commitment process.” Jennifer L. Skeem & John Monahan, *Current Directions in Violence Risk Assessment*, 20 CURRENT DIRECTIONS IN PSYCHOL. SCI. 38, 38 (2011).

⁵² Janus & Prentky, *supra* note 31, at 1444.

⁵³ See Woodworth & Kadane, *supra* note 51, at 236.

⁵⁴ See *id.* (suggesting how risk assessments for violent sexual offenders may be changed to pass *Daubert* challenges); see also Janus & Prentky, *supra* note 31, at 1446 (arguing that actuarial risk assessments in civil commitments are admissible under *Daubert*).

C. Analogy to Evidence-Based Sentencing

The legal and academic discourse regarding future dangerousness predictions in capital cases and civil commitments highlights two competing goals—allowing the introduction of appropriate scientific evidence and giving the fact-finder as much information as possible to render an appropriate decision.⁵⁵ This tension also exists in evidence-based sentencing.

Despite the fact that sentencing hearings do not subject scientific evidence to the same rigorous testing as during trial,⁵⁶ this Note will follow the academic discussion surrounding future dangerousness predictions in capital cases and civil commitments and argue that evidence-based sentencing should pass *Daubert*.⁵⁷

III

DAUBERT'S APPLICATION TO EVIDENCE-BASED SENTENCING GENERALLY AND IN PENNSYLVANIA

“*Daubert* unequivocally endorses ‘empirically validated treatments’ and ‘evidence-based practices’” and protects the jury from giving undue weight to an unreliable and invalid expert testimony.⁵⁸ A *Daubert* analysis requires trial judges to decide “whether the reasoning or methodology underlying the

⁵⁵ See Woodworth & Kadane, *supra* note 51, at 236 (proposing “that the statistical expert’s job is to provide the fact-finder with the means to determine, as accurately and specifically as possible, the probability that this offender will recidivate given what is known”); see also Beecher-Monas & Garcia-Rill, *supra* note 33, at 1901 (noting that “[b]ecause juries are likely to evaluate the future dangerousness of any criminal defendant . . . they should be provided with the most accurate information that can bear on such an assessment. This includes actuarial risk factor studies . . .”).

⁵⁶ One legal commentator characterizes the lack of safeguards around expert testimony at sentencing hearings as “an evidentiary free-for-all.” Beecher-Monas, *supra* note 45, at 357. Nineteen states have ruled that either the states’ evidentiary admissibility standards do not apply to expert testimony based on structured risk assessments or, if they do apply, these structured risk assessments meet the unchallenging standard. Daniel A. Krauss & Nicholas Scurich, *Risk Assessment in the Law: Legal Admissibility, Scientific Validity, and Some Disparities Between Research and Practice*, 31 BEHAV. SCI. L. 215, 220 (2013).

⁵⁷ *But see generally* Pari McGarraugh, Note, *Up or Out: Why “Sufficiently Reliable” Statistical Risk Assessment Is Appropriate at Sentencing and Inappropriate at Parole*, 97 MINN. L. REV. 1079, 1083, 1112 (2013) (arguing for “sufficiently reliable” as an admissibility standard of actuarial risk assessments because of the other due process safeguards that are available to defendants at sentencing).

⁵⁸ David L. Faigman & John Monahan, *Psychological Evidence at the Dawn of the Law’s Scientific Age*, 56 ANN. REV. PSYCHOL. 631, 656 (2005); see Krauss & Scurich, *supra* note 56, at 219.

testimony is scientifically valid and . . . whether that reasoning or methodology properly can be applied to the facts in issue.”⁵⁹

Thus, testimony or evidence⁶⁰ must be reliable *and* relevant to help the trier of fact resolve the factual issue.⁶¹ Legally,⁶² reliability of the scientific technique turns on its scientific validity.⁶³ The relevance requirement under *Daubert* asks whether the information is helpful in deciding the case at hand.⁶⁴ This allows the trial court to analyze “whether this particular expert ha[s] sufficient specialized knowledge to assist the jurors ‘in deciding the particular issues in the case.’”⁶⁵

Despite acknowledging that there is no “definitive checklist or test,”⁶⁶ the United States Supreme Court laid out four factors typically relevant to determining reliability of scientific evidence: 1) “whether a theory or technique . . . can be (and has been) tested”; 2) “whether the theory or technique has been subjected to peer review and publication”;⁶⁷ 3) what the known error rate (or potential rate of error) is and whether standards that control the technique’s operation exist; and 4) whether the technique is generally accepted within a relevant scientific community and to what degree.⁶⁸ This analysis is a “flexible” inquiry focusing on “principles and methodology.”⁶⁹

⁵⁹ *Daubert v. Merrell Dow Pharm., Inc.*, 509 U.S. 579, 592–93 (1993).

⁶⁰ The *Daubert* test also includes expert testimony involving technical and other specialized knowledge. See *Kumho Tire Co. v. Carmichael*, 526 U.S. 137, 141 (1999) (holding “that *Daubert*’s general holding—setting forth the trial judge’s general ‘gatekeeping’ obligation—applies not only to testimony based on ‘scientific’ knowledge, but also to testimony based on ‘technical’ and ‘other specialized’ knowledge”).

⁶¹ See *Daubert*, 509 U.S. at 589 (stating that “under the Rules the trial judge must ensure that any and all scientific testimony or evidence admitted is not only relevant, but reliable”); Janus & Prentky, *supra* note 31, at 1460–61.

⁶² The author would like to draw attention to the fact that in the sciences, validity and reliability are two separate concepts. In such context, validity measures whether the answer given by a tool or test is correct. Reliability measures whether the answer given (regardless of correctness) is consistent. For a discussion of these differences, see DAVID B. ORR, *FUNDAMENTALS OF APPLIED STATISTICS AND SURVEYS* 54–55 (1995).

⁶³ See *Daubert*, 509 U.S. at 590, n.9; Scherr, *supra* note 43, at 9.

⁶⁴ See Scherr, *supra* note 43, at 9. For the sake of brevity, this Note will assume that the evidence-based sentencing is relevant, but readers are encouraged to consider Melissa Hamilton’s article, *Adventures in Risk*, *infra* note 78, specifically her discussion on “Fit,” as well as *infra* Part V (Public Policy Concerns).

⁶⁵ *Kumho Tire Co.*, 526 U.S. at 156.

⁶⁶ *Daubert*, 509 U.S. at 593.

⁶⁷ *Id.* “[S]ubmission to the scrutiny of the scientific community is a component of ‘good science,’ in part because it increases the likelihood that substantive flaws in methodology will be detected.” *Id.*

⁶⁸ *Id.* at 594.

⁶⁹ *Id.* at 594–95.

The following sections analyze evidence-based sentencing under each of the four *Daubert* factors for reliability. Each section draws on the literature surrounding actuarial assessments in general, followed by a closer analysis of the Pennsylvania Commission on Sentencing's proposed tool.

Although more than twenty states have embraced some use of evidence-based sentencing, they do so in very different ways.⁷⁰ In some jurisdictions, state legislatures have passed laws requiring actuarial risk assessments in sentencing procedures,⁷¹ while in other jurisdictions, the judicial branch has endorsed risk assessments not required by the legislature.⁷² To make matters more complicated, in still other jurisdictions, the judiciary and legislative branches have worked in tandem to create an evidence-based sentencing system.⁷³ After deciding to use evidence-based sentencing, states can either use commercial risk assessments, such as LSI-R and COMPAS, or create their own.⁷⁴ Therefore, the discussion below regarding evidence-based sentencing in general will be broad because a survey reviewing such diverse practices is outside of the scope of this Note.⁷⁵

Finally, this Note argues that even if evidence-based sentencing were to pass a *Daubert* analysis, a mechanism such as Federal Rule of Evidence 403 would disallow the testimony. Federal Rule of Evidence 403 gives the Court “‘very broad’ discretion”⁷⁶ to exclude the expert evidence if the probative value

⁷⁰ See Center for Sentencing Initiatives, *State Policies & Legislation: Interactive Map*, NCSC (July 2016), <http://www.ncsc.org/microsites/csi/home/In-the-States/State-Activities/RNA-Map.aspx> [<https://perma.cc/3V7F-6XEK>]; Luis Daniel, *The Dangers of Evidence-Based Sentencing*, GOVLAB (Oct. 31, 2014), <http://thegovlab.org/the-dangers-of-evidence-based-sentencing/> [<https://perma.cc/5X3G-AT6T>].

⁷¹ Some such laws include TENN. CODE ANN. § 41-1-412 (2017); OHIO REV. CODE ANN. § 5120.114 (West 2011); OKLA. STAT. tit. 22 § 22-988.18 (West 2017); LA. STAT. ANN. § 15:327 (2014); COLO. REV. STAT. § 16-11-102 (2017); KY. REV. STAT. ANN. § 532.007 (West 2015); UTAH CODE ANN. § 77-18-1(3) (West 2017).

⁷² Some examples of these jurisdictions are Arkansas, Indiana, and Arizona. See Center for Sentencing Initiatives, *supra* note 70.

⁷³ *Id.*

⁷⁴ For example, Wyoming uses LSI-R, Michigan uses COMPAS, and Ohio has its own risk assessment system (ORAS). *Id.*

⁷⁵ Though a comparison of different categories of current practices would significantly add to current literature, much of this information is inaccessible to the public.

⁷⁶ *United States v. Smithers*, 212 F.3d 306, 322 (6th Cir. 2000) (quoting *United States v. Hawkins*, 969 F.2d 169, 174 (6th Cir. 1992)).

is substantially outweighed by a danger of unfair prejudice⁷⁷ or of confusing the jury.⁷⁸

A. Introduction to the Pennsylvania Commission on Sentencing Risk Assessment Tool

The Pennsylvania Commission on Sentencing started developing a risk assessment tool in 2010.⁷⁹ The Commission developed and validated a risk assessment tool for all Level 3 and Level 4 offenses.⁸⁰ In 2015, the Commission published a report on the development of a new risk assessment scale for all five offense levels in Pennsylvania.⁸¹ The results of this risk assessment tool (based on nine categories of Offense Gravity Scores for the risk of re-offense for any crime) was developed and validated twice using samples of offenders sentenced 1998-2000 and 2004-2006 and published in February 2016.⁸² The Commission used the risk assessment tool created from 1998-2000 development sample, but not yet validated, to analyze the effects of removing sex, age, and county from the model.⁸³ Finally, in February 2016, the Pennsylvania Commission on Sentencing indicated its interest in creating a risk assessment tool to assess the risk of re-offense for a crime against a person.⁸⁴ Although the proposed law implementing the risk assessment instruments indicates that this tool has been created and will be used along with the risk assessment tool assessing the risk of re-offense for any crime,⁸⁵ the Commission

⁷⁷ “Unfair prejudice ‘does not mean the damage to a defendant’s case that results from the legitimate probative force of the evidence; rather it refers to evidence which tends to suggest [a] decision on an improper basis.’” *United States v. Bonds*, 12 F.3d 540, 567 (6th Cir. 1993) (quoting *United States v. Mendez-Ortiz*, 810 F.2d 76, 79 (6th Cir. 1986)).

⁷⁸ See Melissa Hamilton, *Adventures in Risk: Predicting Violent and Sexual Recidivism in Sentencing Law*, 47 ARIZ. ST. L.J. 1, 56–57 (2015) [hereinafter *Adventures in Risk*].

⁷⁹ See PA. COMM’N ON SENT’G, RISK/NEEDS ASSESSMENT PROJECT INTERIM REPORT 1: REVIEW OF FACTORS USED IN RISK ASSESSMENT INSTRUMENTS 2 (2011).

⁸⁰ See PA. COMM’N ON SENT’G, RISK ASSESSMENT PROJECT II INTERIM REPORT 2: VALIDATION OF A RISK ASSESSMENT INSTRUMENT BY OFFENSE GRAVITY SCORE FOR ALL OFFENDERS 2 (2016). [hereinafter *VALIDATION BY OFFENSE GRAVITY SCORE*].

⁸¹ See PA. COMM’N ON SENT’G, RISK/NEEDS ASSESSMENT PROJECT II INTERIM REPORT 1: DEVELOPMENT OF A RISK ASSESSMENT SCALE BY OFFENSE GRAVITY SCORE FOR ALL OFFENDERS 1 (2015).

⁸² See *id.* at 1–4.

⁸³ See PA. COMM’N ON SENT’G, RISK ASSESSMENT PROJECT II SPECIAL REPORT: THE IMPACT OF REMOVING AGE, GENDER, AND COUNTY FROM THE RISK ASSESSMENT SCALE 2 (2015) [hereinafter *IMPACT OF REMOVING DEMOGRAPHIC FACTORS*].

⁸⁴ See *VALIDATION BY OFFENSE GRAVITY SCORE*, *supra* note 80, at 28.

⁸⁵ See PA. COMM’N ON SENT’G, PROPOSALS PUBLISHED IN PENNSYLVANIA BULLETIN (47 PA.B. 1999) § 305.2 (2017) [hereinafter *PROPOSALS PUBLISHED*], <http://>

has yet to release any of the data on the development or validation of such sample in its available reports.⁸⁶

B. Prong 1: Whether Evidence-Based Sentencing Has Been Tested?

The first prong of the *Daubert* analysis requires determining whether evidence-based sentencing can be and has been tested.⁸⁷ This Note will first review the testing of common similar risk-assessment tools (though they are not necessarily only used for sentencing)⁸⁸ to show the difficulty with determining what is being tested and whether it can be tested. This Note will then turn to address the Pennsylvania Commission on Sentencing's risk assessment model.

Before deciding whether evidence-based sentencing can be and has been tested, it is necessary to ask what is being tested. This could include one of two questions: whether the statistical technique used to calculate a risk score has been tested, and whether the outcomes of such a statistical technique have been tested.⁸⁹

As detailed in Part I, evidence-based sentencing relies on determining the predictive factors for a certain developmental

pcs.la.psu.edu/guidelines/proposed-for-public-comment-sentence-risk-assessment-instrument/annex-b/view [https://perma.cc/P5N9-5UV4].

⁸⁶ The reports featured at the time this analysis was done are: Phase I Interim Report 1: Review of Factors Used in Risk Assessment Instruments; Phase I Interim Report 2: Recidivism Study: Initial Recidivism Information; Phase I Interim Report 3: Factors that Predict Recidivism for Various Types of Offenders; Phase I Interim Report 4: Development of Risk Assessment Scale; Phase I Interim Report 5: Developing Categories of Risk; Phase I Interim Report 6: Impact of Risk Assessment Tool for Low Risk Offenders; Phase I Interim Report 7: Validation of Risk Scale; Phase I Interim Report 8: Communicating Risk at Sentencing; Phase I Special Report: The Impact of Juvenile Record on Recidivism Risk; Phase II Interim Report 1: Development of a Risk Assessment Scale by Offense Gravity Score for All Offenders; Phase II Interim Report 2: Validation of Risk Assessment Instrument by OGS for All Offenses; Phase II Special Report: Impact of Removing Demographic Factors; Sentence Risk Assessment Instrument Adopted for Purpose of Public Comment. See PA. COMM'N ON SENT'G, *Overview of the Sentence Risk Assessment Instrument*, <http://pcs.la.psu.edu/publications-and-research/research-and-evaluation-reports/risk-assessment> [https://perma.cc/TY6S-36C8] (follow "Phase I," "Phase II," and "Sentence Risk Assessment Instrument for purposes of public comment" hyperlinks). Subsequent reports and drafts are still being added.

⁸⁷ See *Daubert v. Merrell Down Pharm., Inc.*, 509 U.S. 579, 593 (1993).

⁸⁸ Many tests are wrongly used for sentencing despite specific warnings not to use them. See, e.g., Starr, *supra* note 1, at 812 (citing DEP'T OF CORRECTIVE SERVS., LSI-R TRAINING MANUAL 8 (2002); BERNARD E. HARCOURT, *AGAINST PREDICTION: PROFILING, POLICING, AND PUNISHING IN AN ACTUARIAL AGE 78-84* (2007); Starr, *supra* note 1, at 809 n.11).

⁸⁹ See HANDBOOK OF PSYCHOLOGY IN LEGAL CONTEXTS 375-77 (David Carson & Ray Bull eds., 2d ed. 2003).

or normed sample, attributing appropriate weights for each predictive factor, and creating an estimated probability of the outcome (here, recidivism) occurring for each score or group of scores.⁹⁰ This is often achieved through a logistic regression, a commonly accepted statistical technique.⁹¹ The Pennsylvania Commission on Sentencing used a logistic regression to create its actuarial instrument.⁹²

Although using logistic regression to create predictive instruments is a commonly accepted and tested technique, a closer look at the Pennsylvania Commission on Sentencing's predictive instrument shows that not all logistic regressions are created equally. First, the predictive ability of a test is limited to the development sample.⁹³ For example, the Pennsylvania Commission on Sentencing randomly divided its samples of the offenders convicted between 1998 to 2000 and 2004 to 2006 at all five offense levels into two groups: one development and one validation.⁹⁴ The fact that the models derived from half of these two development groups were validated using the other half of the development groups in the same time frame does little to show its predictive value in the future, which is the purpose of this study. After creating and validating the model by two samples in the same time frame, the Commission should have gone one step further to test the predictive ability of the model on a sample that had not been used to develop the model, such as 2008 to 2010.

For the proposed law implementing these risk assessment tools, recidivism is measured by re-arrests within three years, which does not include out-of-state, federal or foreign charges, or a re-arrest *dismissed* by a minor court.⁹⁵ Although the Pennsylvania Sentencing Commission included technical viola-

⁹⁰ See *Adventures in Risk*, *supra* note 78, at 9. Within the group of score known as "risk bins" (high, medium, low), there is no commonly agreed upon definitions of risk categories, accepted metrics, or normative legal distinctions for such labels. Among different tests, the same person can have a vastly different risk category. See, e.g., *id.* at 21 (describing a study in which sex offenders received disparate categorical labels based on scores from different risk assessment instruments).

⁹¹ Personal communication with Professor John Zipp, University of Akron (Jan. 2, 2017). See GARETH JAMES ET AL., AN INTRODUCTION TO STATISTICAL LEARNING: WITH APPLICATIONS IN R 127 (2013, corrected at 8th printing 2017).

⁹² See PA. COMM'N ON SENT'G, INTERIM REPORT 3: FACTORS THAT PREDICT RECIDIVISM FOR VARIOUS TYPES OF OFFENDERS (PHASE I) 6 (2011) [hereinafter *FACTORS THAT PREDICT RECIDIVISM*].

⁹³ See HANDBOOK OF PSYCHOLOGY IN LEGAL CONTEXTS, *supra* note 89, at 381.

⁹⁴ See VALIDATION BY OFFENSE GRAVITY SCORE, *supra* note 80, at 9.

⁹⁵ PROPOSALS PUBLISHED, *supra* note 85, § 305.1(b)(15).

tions as re-arrests in its development sample,⁹⁶ the final proposed law explicitly excludes the use of the Sentence Risk Assessment Instrument for the sentences imposed as a result of “revocation of probation, intermediate punishment or parole.”⁹⁷ Depending on the quantity and shared characteristics of those re-arrested on a technical violation, the current risk assessment tool is likely to fit less well due to this statutory exclusion that was not accounted for in the creation of the model.

Additionally, arrest itself may be a false proxy for criminality resulting from higher level of policing in certain neighborhoods.⁹⁸ The high incidence of plea bargaining in the U.S. criminal justice system, which may include a significant amount of innocent people who plead guilty even though they could have had their case dismissed by a minor court, compounds concerns that arrest is an inappropriate measurement of recidivism.⁹⁹

In the end, the results from an actuarial risk assessment (usually a number on a scale) do little to determine the qualitative features of the defendant’s punishment (jail, probation, or fine; length of time; and possibility of parole). The Commission attempts to contextualize the recidivism risk by using two separate evidence-based sentencing tools for “Risk of Re-Offense for Any Offense” and the “Risk of Re-Offense for an Offense

⁹⁶ See VALIDATION BY OFFENSE GRAVITY SCORE, *supra* note 80, at 4.

⁹⁷ PROPOSALS PUBLISHED, *supra* note 85, § 305.3(b).

⁹⁸ When “there are more police officers making more arrests in high-crime neighborhoods,” those living in that area “will automatically score higher simply because of where they live.” Therefore, when police are predisposed to arrest minorities without good cause, a prediction based on arrests can falsely interpret arbitrary arrests as a sign of increased criminality among minority groups, which creates a positive feedback loop that reinforces the pattern. See Ryan Briggs, *Should Your Race or ZIP Code Determine Jail Time?*, CITY & STATE PA. (Aug. 2, 2016, 11:39 AM), <http://www.cityandstatepa.com/content/brewing-battle-over-what-factors-will-determine-jail-time> [<https://perma.cc/AMW7-LTZC>]. See also Testimony of the Defender Association of Philadelphia, Pennsylvania Association of Criminal Defense Lawyers, Public Defender Association of Pennsylvania and the Allegheny County Public Defender’s Office Before the Sentencing Commission of Pennsylvania (May 25, 2017), at 3–4, <http://pcs.la.psu.edu/guidelines/sentencing/sentencing-guidelines-and-implementation-manuals/7th-edition-amendment-4-1-1-2018-1/7th-edition-amendment-4-adopted-for-purposes-of-public-comment/testimony/testimony-mark-houldin-and-bradley-s.-bridge-defender-association-of-philadelphia-philadelphia-may-25-2017/view> [<https://perma.cc/G7EW-72RB>] (describing the racial disparity in policing measures in Pennsylvania).

⁹⁹ See Jed S. Rakoff, *Why Innocent People Plead Guilty*, N.Y. REV. OF BOOKS (Nov. 20, 2014), <http://www.nybooks.com/articles/2014/11/20/why-innocent-people-plead-guilty/?insrc=whc> [<https://perma.cc/S6KS-GGKB>].

Against a Person”;¹⁰⁰ however, the contours and relative gravity of these crimes are hazy. A robbery of a motor vehicle is considered an “Offense Against a Person,”¹⁰¹ but operating a meth lab is not.¹⁰² Similarly, gravity between and among these types of crimes is not intuitive; for example, both fortunetelling for money and operating a meth lab are considered “Any Offense.”¹⁰³ Thus, the Pennsylvania Commission on Sentencing’s predictions tell us the probability a person will be re-arrested within the next three years for a wide range of crimes based on a tool that has not been validated on a future sample but instead only on a sample in the same time frame as the model’s development sample.

Although choosing the right sample upon which to create a predictive model through logistic regression is crucial, the Pennsylvania Commission on Sentencing’s methods also demonstrate how the selection of the model’s variables is important. When selecting the factors to include in its model, Pennsylvania’s risk assessment tool only took the factors found to be significant predictors of recidivism.¹⁰⁴ Although race and county were statistically significant factors,¹⁰⁵ the final Pennsylvania Commission on Sentencing risk assessment tool does not include race and county as factors.¹⁰⁶ The Commission explains that race and county are “statistically controlled for in the analyses, which means that the effects of the other factors are included only after eliminating the effects of race and county.”¹⁰⁷ The Commission never explains why it keeps race and county in the model to create the risk assessment tool even when it knows that it is going to exclude it from the risk scale. Thus, it is preferable that the logistic regression model be re-run without race and county rather than controlling for it post hoc. By including both race and county in the model, the Commission could be ignoring other factors that could have been

¹⁰⁰ See PROPOSALS PUBLISHED, *supra* note 85, § 305.2(b)(2).

¹⁰¹ See PROPOSALS PUBLISHED, *supra* note 85, § 305.1(b)(9) (citing 42 PA. CONS. STAT. § 9714(g) (2015)).

¹⁰² Given that the definition of § 305.1(b)(9) in PROPOSALS PUBLISHED, *supra* note 85, is exclusive, operating a meth lab would not be considered an “Offense Against a Person” as it does not fall within any of the three categories of crimes defining an “Offense Against a Person.” See 18 PA. CONS. STAT. § 7508.2 (2015) (operating a meth lab).

¹⁰³ Because neither is included in the definition of “Offense Against a Person,” they both fall into the residual category of “Any Offense.” See 18 PA. CONS. STAT. § 7104 (2015) (fortunetelling).

¹⁰⁴ See VALIDATION BY OFFENSE GRAVITY SCORE, *supra* note 80, at 12.

¹⁰⁵ See *id.* at 51–58 (“Final Logistic Models by Offense Gravity Score”).

¹⁰⁶ See PROPOSALS PUBLISHED, *supra* note 85, § 305.1(b)(18).

¹⁰⁷ VALIDATION BY OFFENSE GRAVITY SCORE, *supra* note 80, at 12 n.8.

significant in the development model had race and county not been initially included.¹⁰⁸

The outcomes of such predictive instruments have been typically measured using area under the curve (AUC) analysis to measure the discrimination of the test; that is, how well the predictive test can differentiate those who experienced the outcome of interest (labeled as likely to recidivate and did recidivate) versus those who did not (labeled as likely to recidivate but did not recidivate).¹⁰⁹ Melissa Hamilton notes that many researchers misconstrue the meaning of the area under the curve when evaluating risk assessment tools. The AUC analysis provides “‘the probability that a randomly selected individual who committed an [act of recidivism] . . . received a higher risk classification than a randomly selected individual who did not’ reoffend.”¹¹⁰ This does not tell us “information on the accuracy of any individual prediction,”¹¹¹ “the probability that individuals are scored correctly . . . [, or] the potential that a person assessed with a high test score will eventually become a recidivist.”¹¹²

Typically, AUC values in popular risk assessments for violent and sexual recidivism tools range from 70-75%, meaning that “these risk instruments have been able to classify violent and sexual recidivists at higher levels of risk than non-recidivists about 70 to 75% of the time.”¹¹³ Although these values

¹⁰⁸ Personal Communication with Professor John Zipp (Jan. 2, 2017). See also generally *Adventures in Risk*, *supra* note 78, at 20–21 (discussing how “risk tools typically include a relatively small number of variables, thereby omitting a plethora of potential explanatory or correlative factors”).

¹⁰⁹ See *Adventures in Risk*, *supra* note 78, at 25.

¹¹⁰ *Id.* at 25 (alteration in the original) (quoting Jay P. Singh et al., *Measurement of Predictive Validity in Violence Risk Assessment Studies: A Second-Order Systematic Review*, 31 BEHAV. SCI. & L. 55, 64 (2013)). “An AUC of .90, as an illustration, means that if one randomly chooses a recidivist and a non-recidivist, the recidivist’s actuarial score would be higher than the non-recidivist’s score about 90% of the time.” *Id.* (citing Christopher T. Lowenkamp et al., *The Federal Post Conviction Risk Assessment (PCRA): A Construction and Validation Study*, 10 PSYCHOL. SERVS. 87, 92 n.11 (2013)). More problematic is that “a scale can achieve a high rating for discrimination even when the average predicted risk of violent re-offense is significantly different than the actual percentage of violent recidivists.” *Id.* at 24 (citing Nancy R. Cook, *Use and Misuse of the Receiver Operating Characteristic Curve in Risk Prediction*, 115 CIRCULATION 928, 928 (2007)).

¹¹¹ *Id.* at 27 (citing Thomas Nilsson et al., *The Precarious Practice of Forensic Psychiatric Risk Assessments*, 32 INT’L J.L. & PSYCHIATRY 400, 402 (2009)).

¹¹² *Id.* at 26 (citing Cook, *supra* note 110, at 928).

¹¹³ *Id.* Some experts contend that this common rate is due to similar factors being used by recidivism risk tools and that technologies have possibly reached the natural limit for predicting human behavior. *Id.* at 28 (citing John Monahan & Jennifer L. Skeem, *Risk Redux: The Resurgence of Risk Assessment in Criminal Sanctioning*, 26 FED. SENT’G REP. 158, 158 (2014)).

perform statistically better than chance (50%), many researchers equate these values to moderate or large effect sizes.¹¹⁴ However, in statistics, there is no consensus as to which AUC scores represent small, moderate, or large effect sizes.¹¹⁵ Without such a consensus, researchers could tout that a certain predictive instrument effect size is large enough to be used in the courts without having to define what “large” is across the literature or describe the specific definition of a tool’s AUC.¹¹⁶

Finally, the Commission confines the use of this tool in a way that may not fully achieve the benefits of evidence-based sentencing, particularly reducing incarceration rates. For offenders falling into the typical risk category (the middle 68%), the Commission makes no additional recommendation; therefore these defendants are sentenced on the information available.¹¹⁷ For low- or high-risk offenders, the Commission recommends a separate pre-sentencing report to collect more information before sentencing.¹¹⁸ Whether more closely analyzing the remaining 32% of low- or high-risk offenders’ situations will actually result in different, more effective incarceration is an open question. Although using logistic regressions to predict individual outcomes as a general technique can be and has been tested, the Pennsylvania Commission on Sentencing fails to justify why its sample, validation method, and variables differ from what one would expect in a similar scenario.

C. Prong 2: Whether Evidence-Based Sentencing Has Been Subjected to Peer Review and Publication?

Judges typically pay close attention to publication and peer review, even more so than the testing or error rates.¹¹⁹ Judges may focus on journal names as a proxy for quality scholarship, but it is important to assess the content of the published studies.¹²⁰ The importance of peer review is to de-

¹¹⁴ *Id.* at 26 (citing R. Karl Hanson & David Thornton, *Improving Risk Assessments for Sex Offenders: A Comparison of Three Actuarial Scales*, 24 LAW & HUM. BEHAV. 119, 129 (2000)).

¹¹⁵ *Id.*

¹¹⁶ Hamilton suggests that focusing on calibration is a better benchmark for evaluating an instrument’s predictive ability. *Id.* at 28–29.

¹¹⁷ See PROPOSALS PUBLISHED, *supra* note 85, §§ 305.1(b)(17), 305.5(b).

¹¹⁸ See PROPOSALS PUBLISHED, *supra* note 85, § 305.5(c).

¹¹⁹ Jay P. Kesan, *An Autopsy of Scientific Evidence in a Post-Daubert World*, 84 GEO. L.J. 1985, 2029 (1996).

¹²⁰ *Id.* at 2029; *Adventures in Risk*, *supra* note 78, at 57–61. See generally MICHAEL J. SAKS ET AL., ANNOTATED REFERENCE MANUAL ON SCIENTIFIC EVIDENCE

tect the design flaws of studies; any conclusion is only as good as the methods used to attain it.¹²¹

Generally, evidence-based sentencing has received policy-based and learning-oriented discussions.¹²² Peer reviews specific to the technique and efficacy of actuarial risk assessments, in general, are commonly limited to predictions of violence,¹²³ not recidivism. With its growing use, however, there is an increasing presence of academic literature surrounding risk assessment in sentencing.¹²⁴

Several factors impede the ability of the scientific community to peer review methods specific to validity. First, the proprietary nature of commercial tests creates conditions in which the company producing the test is validating its own studies or is involved in another researcher's validation.¹²⁵ This self-review does not count as peer review, and those reviewing it on behalf of the company should disclose this potential conflict of interest. Additionally, limited data from states where evidence-based sentencing has been used for a long time, such as Virginia, make it difficult to evaluate both the methods and efficacy.¹²⁶

Not only does the proprietary nature of commercial tools impede the ability of scientists to conduct peer review on the

107–08 (2d ed. 2004) (discussing several scenarios where a position's illusion of prestige masks flaws in their methods).

¹²¹ See Kesan, *supra* note 119, at 2029–30.

¹²² See, e.g., PAMELA M. CASEY ET AL., NAT'L CTR. FOR STATE COURTS, OFFENDER RISK & NEEDS ASSESSMENT INSTRUMENTS: A PRIMER FOR COURTS 2 (2014), [hereinafter OFFENDER RISK & NEEDS ASSESSMENT INSTRUMENTS] http://www.ncsc.org/~media/Microsites/Files/CSI/BJA%20RNA%20Final%20Report_Combined%20Files%208-22-14.ashx [<https://perma.cc/MXU7-X5NZ>] (“This Primer is a resource to help judges and others involved in sentencing understand and make knowledgeable decisions about the value and use of a [risk and needs] assessment.”).

¹²³ See Jay P. Singh et al., *Measurement of Predictive Validity in Violence Risk Assessment Studies: A Second-Order Systematic Review*, 31 BEHAV. SCI. L. 55, 55 (2013) (“[T]hese [risk assessment] instruments are designed to aid in the assessment of risk for antisocial behavior, most commonly general violence, sexual violence, and criminal offending.”); Min Yang et al., *The Efficacy of Violence Prediction: A Meta-Analytic Comparison of Nine Risk Assessment Tools*, 136 PSYCHOL. L. BULL. 740, 741 (2010) (“In most cases, actuarial tools are designed by combining empirically or theoretically derived constructs that are predictive of violence or antisocial activities to guide the forecasting of future antisocial or violent acts.”).

¹²⁴ Skimming through the following sources helps give an idea of the amount of publications. Letter from Brian D. Johnson, Professor at Univ. of Md., James P. Lynch, Professor at Univ. of Md., & Rebecca Richardson, Doctoral Student, to Conn. Sentencing Comm'n (Dec. 1, 2015), <http://www.ct.gov/ctsc/lib/ctsc/2016-01DRAFT.pdf> [<https://perma.cc/52RG-TP8P>]; OFFENDER RISK & NEEDS ASSESSMENT INSTRUMENTS, *supra* note 122, at 32.

¹²⁵ See Starr, *supra* note 8 (noting that risk prediction instruments are frequently not very transparent).

¹²⁶ *The New Science of Sentencing*, *supra* note 24.

tool, it also interferes with due process considerations. In March 2017, the United States Supreme Court indicated that it was interested in whether a defendant's right to due process was violated by a judge's consideration of a COMPAS report that the defendant was unable to inspect or challenge by requesting that the federal government file an *amicus* brief arguing whether it should hear the case.¹²⁷

Although there is a growing body of peer review literature for evidence-based sentencing, Hamilton's critique of incorrect statistical interpretation warns of the danger of reading too quickly.¹²⁸ Under *Daubert*, courts should consider which test is being presented when looking for specific peer review, as not all research on risk assessment will be applicable. Although the Pennsylvania Commission on Sentencing has said that it will seek an external validation once the development process is completed, if it has occurred, it is not yet available, and the risk assessment program is set to begin on July 1, 2018 without any restrictions subjecting the risk assessment tools to external peer review.¹²⁹

D. Prong 3: What Is the Known (or Potential) Error Rate and Are There Standards that Control Evidence-Based Sentencing?

The error rates of evidence-based sentencing are crucial to the *Daubert* analysis and directly hinge on how much error a society is willing to tolerate in exchange for the purported safety benefits of correctly identifying a defendant with a high-risk of recidivism.¹³⁰ Evidence-based sentencing, by churning

¹²⁷ Adam Liptak, *Sent to Prison by a Software Program's Secret Algorithms*, N.Y. TIMES (May 1, 2017), <https://www.nytimes.com/2017/05/01/us/politics/sent-to-prison-by-a-software-programs-secret-algorithms.html> [perma.cc/VR2Z-J422].

¹²⁸ *Adventures in Risk*, *supra* note 78, at 2–3.

¹²⁹ The fact that Pennsylvania is seeking an external peer review is available through a source created by the former Chair of the Pennsylvania Commission on Sentencing. No other information on the status of this external validation was posted on the Commission on Sentencing's risk assessment website. Steven L. Chanenson & Jordan M. Hyatt, *The Use of Risk Assessment at Sentencing: Implications for Research and Policy* 9 (Villanova U. Charles Widger Sch. of L. Working Paper Series, 2016) <http://digitalcommons.law.villanova.edu/cgi/viewcontent.cgi?article=1201&context=wps> [https://perma.cc/G855-3ZDP].

¹³⁰ Some try to draw false negative parallels with a comparison between medicine and law, but law has different standards and goals than medicine. See *Adventures in Risk*, *supra* note 78, at 55. *But see*, Aaron E. Carroll, *What We Mean When We Say Evidence-Based Medicine*, N.Y. TIMES (Dec. 27, 2017), <https://www.nytimes.com/2017/12/27/upshot/what-we-mean-when-we-say-evidence-based-medicine.html> [https://perma.cc/988D-KHK6] (explaining how even “evi-

out a single number on a scale interpreted as an indication of a defendant's future recidivism rate, has the potential to mask the risk of error with the allure of the scientific method.¹³¹

Each risk assessment tool had its own calculable error rate, which means that each tool should be individually subject to *Daubert*. However, to provide some context to the general efficacy of risk assessments, a meta-analysis found that nine different violence risk assessment tools had comparable, moderate levels of predictive efficacy and recommended that they not be used as the primary means for a criminal justice decision.¹³²

Another problem that aggravates errors is the failure to conduct local validations of commercial methods in approximately 70 percent of the jurisdictions that use assessment instruments, which can result in inaccurate classifications.¹³³ Local validations are important because of the difference between jurisdictions as to their “laws, policies, or sentencing guidelines (e.g., management of felony offenses, eligibility for probation, risk of recidivism defined locally as re-arrest vs. re-conviction) may create a unique set of circumstances and constraints that can reduce the effectiveness of a tool created elsewhere.”¹³⁴

For Pennsylvania, the error rate depending upon the distribution of the risk categories led to an overall prediction accuracy ranging from 49 to 65 percent for the tools developed for the Level 3 and Level 4 samples.¹³⁵

The lower the cutoff point for determining high risk, the better the prediction accuracy for high risk . . . but the worse the prediction accuracy for low risk Alternatively, the higher the cutoff point for determining high risk, the worse the prediction accuracy for high risk, but the better the prediction accuracy for low risk.¹³⁶

dence-based medicine” is facing similar definitional, practical, and policy concerns).

¹³¹ See *id.* at 56–57.

¹³² See Yang et al., *supra* note 123, at 761.

¹³³ PAMELA M. CASEY ET AL., NAT'L CTR. FOR STATE COURTS, USING OFFENDER RISK AND NEEDS ASSESSMENT INFORMATION AT SENTENCING: GUIDANCE FOR COURTS FROM A NATIONAL WORKING GROUP 29–31 (2011), <http://www.ncsc.org/~media/Microsites/Files/CSI/RNA%20Guide%20Final.ashx> [perma.cc/XU2U-TXEB].

¹³⁴ *Id.* at 30 (emphasis omitted).

¹³⁵ PA. COMM'N ON SENT'G, RISK/NEEDS ASSESSMENT PROJECT INTERIM REPORT 7: VALIDATION OF RISK SCALE 2, 4 (2013).

¹³⁶ *Id.* at 4.

The Commission does not provide a similarly user-friendly analysis for its current risk assessment tool that includes all of the levels of offense.¹³⁷

E. Prong 4: Whether Evidence-Based Sentencing Is Generally Accepted Within a Relevant Scientific Community and to What Degree?

The purpose of the general acceptance prong of *Daubert* is to uncover whether there is general agreement that the scientific theory is “not based on a novel theory or procedure that is ‘mere speculation or conjecture.’”¹³⁸ A theory that does not have the acceptance of most of the pertinent scientific community or is disfavored by a substantial part of the scientific community will not be generally accepted.¹³⁹ Although absolute certainty is not necessary for a finding of general acceptance,¹⁴⁰ it is crucial to first determine the scope of the relevant scientific community in order to establish whether evidence-based sentencing is generally accepted.¹⁴¹ The criteria used to decide who is part of the relevant scientific community must be broad enough to reach a spectrum of scientists who are not entirely financially, professionally, and ideologically invested in the field.¹⁴² However, in the search to find a broad enough scientific community, the factfinder must still limit the community to those who truly understand the research basis for the work.¹⁴³ Though judges, parole officers, pre-sentence investigators, corrections officers, and other actors in the criminal justice system are exposed to and use a variety of actuarial risk

¹³⁷ VALIDATION BY OFFENSE GRAVITY SCORE, *supra* note 80, at 22.

¹³⁸ *United States v. Bonds*, 12 F.3d 540, 562 (6th Cir. 1993) (quoting *United States v. Brown*, 557 F.2d 541, 559 (6th Cir. 1977)).

¹³⁹ *Id.* at 562 (citing *Novak v. United States*, 865 F.2d 718, 725 (6th Cir. 1989)).

¹⁴⁰ *Id.*

¹⁴¹ David L. Faigman et al., *Group to Individual (G2i) Inference in Scientific Expert Testimony*, 81 U. CHI. L. REV. 417, 460, 463 (2014) [hereinafter *G2i*] (“The key difficulty [] involves identifying what field should be selected for review. Very often fields are defined by self-interested practitioners or established guilds. For example, if a court asks experts in the areas of polygraphs, bite marks, bullet lead, or hair identification about general acceptance, it would likely hear a chorus of consensus, though each of these areas of claimed expertise has been thoroughly discredited.”).

¹⁴² *Id.* at 461 (“Professional overinvestment might include a researcher who is a leading figure in the field but whose life’s work depends on acceptance of the expertise. Ideological investment might include a researcher whose judgment about the validity of an empirical framework will be influenced by its ability to further a desired policy outcome.”).

¹⁴³ *Id.* at 461–62.

assessments in making decisions, they are not the relevant scientific community.¹⁴⁴

For evidence-based sentencing, the relevant scientific community could extend to statisticians and social scientists, but should exclude, or limit the weight given to, statisticians or researchers who are engaged for financial reasons or those who engage with actuarial risk assessments for purely policy reasons. The first limitation on financial motivations should discount evidence from researchers who have created actuarial risk assessments tools for commercial uses as well as those who have tested these tools using funds from a commercial provider. The second limitation of ideological investment may be more difficult to outline, given that many researchers creating specific tools are doing so with the hopes that it will reduce incarceration rates.

Courts have typically noted that general acceptance is needed with respect to both the theory and the methodology of the technique.¹⁴⁵ Evidence-based sentencing uses group statistics to create an individual prediction. This group to individual prediction (G2i) is common throughout many different fields.¹⁴⁶ The methodology of evidence-based sentencing involves the development of the sample to create the predictive test, the variables chosen, and the interpretation of results typically through logistic regression,¹⁴⁷ which, given the amount of diverse predictive tools,¹⁴⁸ has been generally accepted. At its broadest, logistic regression resulting in individual prediction of an individual is generally accepted as an appropriate technique when analyzing a dichotomous dependent variable.¹⁴⁹

¹⁴⁴ See *Williamson v. Reynolds*, 904 F. Supp. 1529, 1558 (E.D. Okla. 1995) (noting that general acceptance of hair analysis should not be based off “hair experts who are generally technicians testifying for the prosecution,” but rather “scientists who can objectively evaluate such evidence”).

¹⁴⁵ *United States v. Bonds*, 12 F.3d 540, 562 (6th Cir. 1993).

¹⁴⁶ See *G2i*, *supra* note 141, at 421–22 (describing G2i examples such as in medicine, where research regarding at what age women should begin having annual mammograms, provides an empirical framework to help doctors make individual decisions and help manage the risks of breast cancer in individual patients, and in meteorology, where group data models the trajectory and severity of storms to help determine whether to evacuate a community due to the threat of a particular storm).

¹⁴⁷ See, e.g., *Chanenson & Hyatt*, *supra* note 129, at 9 (explaining the methodology of Pennsylvania’s evidence-based sentencing).

¹⁴⁸ PA. COMM’N ON SENT’G, RISK/NEEDS ASSESSMENT PROJECT INTERIM REPORT 1: REVIEW OF FACTORS USED IN RISK ASSESSMENT INSTRUMENTS app. C (2011).

¹⁴⁹ Personal communication with John Zipp (Jan. 2, 2017). See JEFFERY T. ULMER, SOCIAL WORLDS OF SENTENCING: COURT COMMUNITIES UNDER SENTENCING GUIDELINES 43 (Austin T. Turk ed., 1997) (“Although other techniques for modeling

Though not a legal question, what remains to be generally accepted is the question: Who should decide the *policy* rationales for evidence-based sentencing, and which goals of evidence-based sentencing to promote?

IV

DAUBERT'S INTERACTION WITH FEDERAL RULE OF EVIDENCE 403

Interestingly, most courts tend not to consider the *Daubert* factors in a step-by-step analysis.¹⁵⁰ Instead, one study found that the best predictor of whether evidence would be admissible was, *inter alia*, whether the evidence was prejudicial. This section focuses on prejudicial information through the lens of Federal Rule of Evidence 403.

Federal Rule of Evidence 403 acts as a final check on the *Daubert* analysis. Even if a scientific technique passes the *Daubert* analysis of reliability and fit, Federal Rule of Evidence 403 may disallow testimony if the probative value is substantially outweighed for many reasons.¹⁵¹ With regards to evidence-based sentencing, the danger of creating unfair prejudice or confusing the jury would substantially outweigh its probative value. Evidence-based sentencing creates a risk of unfair prejudice by using immutable factors outside of the defendant's control.¹⁵² It also confuses the jury members, who do not engage in appropriate critical inquiries about its scientific validity.¹⁵³

Although a typical evidence-based sentencing analysis contains many questionable factors unrelated to the crime,¹⁵⁴ the subsequent analysis will focus on sex and age. These two categories are analyzed for two reasons. First, the Pennsylvania Commission on Sentencing analyzed the efficacy of its

dichotomous dependent variables exist . . . logistic regression is widely held to be the most appropriate procedure.”).

¹⁵⁰ A. Leah Vickers, *Daubert, Critique and Interpretation: What Empirical Studies Tell Us About the Application of Daubert*, 40 U.S.F. L. REV. 109, 133 (2005).

¹⁵¹ FED. R. EVID. 403 (“The court may exclude relevant evidence if its probative value is substantially outweighed by a danger of one or more of the following: unfair prejudice, confusing the issues, misleading the jury, undue delay, wasting time, or needlessly presenting cumulative evidence.”).

¹⁵² See IMPACT OF REMOVING DEMOGRAPHIC FACTORS, *supra* note 83, at 1. See also Starr, *supra* note 1, at 829 (“[A] generalization about a behavioral tendency like criminal recidivism is simply not comparable to a physical difference . . . [and] certain kinds of generalizations (including those concerning gender) are particularly socially harmful or expressively invidious, even if they have statistical support.”).

¹⁵³ *Adventures in Risk*, *supra* note 78, at 56–57.

¹⁵⁴ Starr, *supra* note 1, at 805.

evidence-based sentencing model when demographic factors were removed, thus showing the relative “costs” in terms of accuracy of not including such questionable factors.¹⁵⁵ Because the Commission only conducted this analysis on the risk assessment tool based on the 1998 to 2000 development model (and not the tool used for the final proposed risk assessment),¹⁵⁶ the figures discussed below may not correspond exactly to the current proposed risk assessment tools. Regardless, the general concepts, critiques, and concerns remain the same.

Additionally, each of these categories touches on a constitutional issue or an issue of concern for the court. The inclusion of sex in evidence-based sentencing creates an equal protection issue subject to heightened scrutiny.¹⁵⁷ Age, although not a protected class and subject to only rational basis review, has been considered by the Supreme Court, particularly when focusing on the sentencing of juvenile offenders.¹⁵⁸

It should be noted that the Pennsylvania Commission on Sentencing’s review of these factors was prompted by the constitutional, ethical, and fairness issues raised by Sonja Starr in *Sentencing by Numbers* and *Evidence-Based Sentencing and the Scientific Rationalization of Discrimination*.¹⁵⁹ Despite the concerns and their moderate accuracy, the Pennsylvania Commission on Sentencing has recommended keeping these factors in current and future risk assessment models.¹⁶⁰

A. Sex

Sex is commonly used in evidence-based sentencing and results in harsher sentences for men.¹⁶¹ Within the Pennsylvania Commission on Sentencing’s eight Offender Gravity

¹⁵⁵ IMPACT OF REMOVING DEMOGRAPHIC FACTORS, *supra* note 83, at 1.

¹⁵⁶ VALIDATION BY OFFENSE GRAVITY SCORE, *supra* note 80, at 21–22.

¹⁵⁷ *Craig v. Boren*, 429 U.S. 190, 197 (1976) (“[C]lassifications by gender must serve important governmental objectives and must be substantially related to achievement of those objectives.”).

¹⁵⁸ See *Miller v. Alabama*, 567 U.S. 460, 489 (2012) (holding that mandatory life without parole for a juvenile convicted of murder is unconstitutional); *Graham v. Florida*, 560 U.S. 48, 79–80 (2010) (ruling that life imprisonment without parole for non-homicide offenses committed by juveniles is unconstitutional); *Roper v. Simmons*, 543 U.S. 551, 570–71 (2005) (holding that the death penalty for juveniles is unconstitutional due to their lack of culpability relative to adults).

¹⁵⁹ IMPACT OF REMOVING DEMOGRAPHIC FACTORS, *supra* note 83, at 1.

¹⁶⁰ *Id.* at 1–2.

¹⁶¹ See *Massie*, *supra* note 25, at 522.

Categories¹⁶² that use sex in their risk assessment, sex accounts for 1-5% of the risk model's prediction of recidivism.¹⁶³ In evaluating whether sex could be removed from the risk assessment model, the Pennsylvania Commission on Sentencing concluded, "[E]ven though removing [sex] . . . would impact the utility of the scale to a much lesser extent [than age],¹⁶⁴ since [it] do[es] provide some increase in the recidivism prediction, it would be worth keeping [it] in the scale as well."¹⁶⁵ Removing sex would result in females appearing to have higher recidivism rates, and males appearing to have lower recidivism rates, as opposed to how these statistics appear in the original scale.

In its most recent validation of Risk Assessment Instrument by Offense Gravity Score [OGS] for All Offenders for the 2004-2006 development sample, male defendants all scored one point higher than female defendants.¹⁶⁶ Given that the total risk scores range from eight to nineteen depending upon the defendant's OGS category, sex carries a substantial weight on the severity of one's sentence.

Despite the fact that the Federal Sentencing Guidelines explicitly forbid sex as a sentencing consideration and that case law regarding the use of sex in non-evidence-based sentencing proceedings is clear,¹⁶⁷ the revised Model Penal Code encourages the inclusion of sex as a factor while offering no equal protection analysis as to why it would be acceptable.¹⁶⁸

¹⁶² The current proposed risk assessment tool for "Any Crime" and "Crime Against a Person" take sex into consideration for seven and six Offender Gravity Categories, respectively. PROPOSALS PUBLISHED, *supra* note 85, § 305.7(a)-(b).

¹⁶³ IMPACT OF REMOVING DEMOGRAPHIC FACTORS, *supra* note 83, at 5.

¹⁶⁴ See *infra* subpart IV.B discussion on Age.

¹⁶⁵ IMPACT OF REMOVING DEMOGRAPHIC FACTORS, *supra* note 83, at 1.

¹⁶⁶ VALIDATION BY OFFENSE GRAVITY SCORE, *supra* note 80, at 19.

¹⁶⁷ U.S. SENTENCING GUIDELINES MANUAL § 5H1.10 (U.S. SENTENCING COMM'N 2012). The few times this issue has arisen, courts have ruled that basing sentences on sex is unconstitutional. See Starr, *supra* note 1, at 824 (citing *United States v. Maples*, 501 F.2d 985, 987 (4th Cir. 1974); *Williams v. Currie*, 103 F. Supp. 2d 858, 868 (M.D.N.C. 2000); Carissa Byrne Hessick, *Race and Gender as Explicit Sentencing Factors*, 14 J. GENDER RACE & JUST. 127, 137, n.68 (2010)).

¹⁶⁸ MODEL PENAL CODE: SENTENCING § 6B.09 reporter's note at 62 (AM. LAW INST., Tentative Draft No. 2, 2011); Starr, *supra* note 1, at 824 (noting that "the MPC drafters . . . offer no commentary defending [the inclusion of sex] on constitutional grounds, as though its constitutionality is self-evident"). The Model Penal Code does note that referring to one's race or ethnicity during sentencing would be unconstitutional. Despite the omission of a seemingly important distinction between the difference in constitutionality of race and sex, the Code was passed in 2011 and remains the same in the Third Tentative Draft. MODEL PENAL CODE: SENTENCING xii (AM. LAW INST., Tentative Draft No. 3, 2014).

This is troubling, because an equal protection analysis would suggest otherwise.¹⁶⁹ Under intermediate scrutiny, distinctions made on the basis of sex must serve “important governmental objectives and [] the discriminatory means employed [must be] substantially related to the achievement of those objectives.”¹⁷⁰ This discrimination cannot be based on “overbroad generalizations” about the differences between men and women.¹⁷¹ Even when data support an even stronger predictive empirical relationship than evidence-based sentencing, the Supreme Court has typically found the discriminatory means to be unconstitutional.¹⁷² This consistent rejection of a sex-based distinction calls into question its use during the sentencing procedure.

B. Age

Age is the most important demographic factor predicting arrest.¹⁷³ In fact, the Pennsylvania Commission on Sentencing found that for each year of increase in age, the likelihood of

¹⁶⁹ See Testimony of the Defender Association of Philadelphia, *supra* note 98, at 5 (focusing on the Pennsylvania Risk Assessment Instrument and arguing that it is “unconstitutional to attach demerit points simply because of one’s gender”). Some scholars argue that considering a variable subject to intermediate or heightened scrutiny along with other variables could be acceptable based on analysis similar to higher education case law regarding affirmative action. These arguments ignore the fact that many other immutable factors are used in this process, that the stakes are not simply one’s educational opportunity but one’s liberty, and that much of the higher education case law focuses on remedying historical oppression through measures such as affirmative action. See Oleson, *supra* note 8, at 1338.

¹⁷⁰ *United States v. Virginia*, 518 U.S. 515, 533 (1996) (internal quotation marks omitted).

¹⁷¹ *Id.* at 532–34.

¹⁷² See Starr, *supra* note 1, at 825–26 (noting that cases, such as *Craig v. Boren*, *Weinberger v. Wiesenfeld*, and *Frontiero v. Richardson* showed that strong statistical data between different sexes did not convince the Court of the constitutionality of classifying by sex based on stereotyping). In *Craig*, men were ten times more likely to be arrested for drunk driving than women. Despite this finding, the Court ruled that the drinking law’s distinction based on sex was unconstitutional. *Craig v. Boren*, 429 U.S. 190, 197 (1976). In *Weinberger v. Frontiero*, women were much more likely to be financially dependent on their husbands (in the 1970s) than the other way around. Despite this finding, the Court ruled that a financially-dependent man was entitled to his wife’s benefits. See *Weinberger v. Wiesenfeld*, 420 U.S. 636, 645 (1975); *Frontiero v. Richardson*, 411 U.S. 677, 690–91 (1973).

¹⁷³ “Age was the most important demographic factor predicting arrest.” IMPACT OF REMOVING DEMOGRAPHIC FACTORS, *supra* note 83, at 1. It is unclear whether this refers to a first or second arrest.

recidivism decreases by roughly 5%.¹⁷⁴ Within the Pennsylvania Commission on Sentencing's nine OGS, age accounts for 21–30% of the risk model's prediction of recidivism.¹⁷⁵ In evaluating whether it could remove this demographic factor from the analysis, the Pennsylvania Commission on Sentencing concluded, "[T]he removal of that factor would have the most impact on recidivism prediction and scale accuracy."¹⁷⁶ If age were removed, older offenders would appear to have higher recidivism rates, whereas younger offenders would appear to have lower recidivism rates as compared to the original scale.¹⁷⁷

In its most recent validation of Risk Assessment Instrument by OGS for All Offenders for the 2004–2006 development sample, defendants under the age of 21 or younger scored three to five points higher than defendants aged 49 or older.¹⁷⁸ Given that the total risk scores range from eight to nineteen depending upon the defendant's OGS category, age carries a significant weight on the severity of one's sentence. The younger a defendant is, the longer their sentence will be.

This result contradicts the Supreme Court's recent jurisprudence surrounding adolescents. Beginning in 2005, the Supreme Court recognized that adolescents are categorically less culpable than adults.¹⁷⁹ Much of the Supreme Court's logic in subsequent decisions involving juveniles was based off of *amici* briefs showing that "developments in psychology and brain science continue to show fundamental differences between juvenile and adult minds."¹⁸⁰ This scientific evidence notes, "Risk taking [including criminal activity] . . . is so pervasive that it 'is statistically aberrant to refrain from such behav-

¹⁷⁴ FACTORS THAT PREDICT RECIDIVISM, *supra* note 92, at 9. The author finds it highly unlikely that the researchers intended to convey a linear effect, particularly given that age is grouped into categories.

¹⁷⁵ IMPACT OF REMOVING DEMOGRAPHIC FACTORS, *supra* note 83, at 5.

¹⁷⁶ *Id.* at 19.

¹⁷⁷ *Id.* at 9.

¹⁷⁸ VALIDATION OF OFFENSE GRAVITY SCORE, *supra* note 80, at 19.

¹⁷⁹ "[P]sychology and brain science continue to show fundamental differences between juvenile and adult minds" making their actions "less likely to be evidence of 'irretrievably depraved character' than are the actions of adults." *Miller v. Alabama*, 567 U.S. 460, 471–72, 490 (2012) (quoting *Roper v. Simmons*, 543 U.S. 551, 570 (2005)).

¹⁸⁰ *Graham v. Florida*, 560 U.S. 48, 68 (2010).

ior during adolescence.”¹⁸¹ Both the Supreme Court and *amici* focus on young persons’ potential for rehabilitation.¹⁸²

Because the Pennsylvania Commission on Sentencing insists on including age in its risk assessment tool, age creates an unfair prejudice. Even though younger people are arguably both legally and scientifically less culpable and more amenable to rehabilitation than their older counterparts, they will likely spend more time in prison.

V

PUBLIC POLICY CONCERNS

Proponents of evidence-based sentencing laud it as the solution for discrimination, overcrowding, and budget shortages.¹⁸³ The discussion on its scientific strength above should serve as a warning that evidence-based sentencing may not be the most effective solution.

There are four major penological theories: retribution, deterrence, incapacitation, and rehabilitation.¹⁸⁴ Over time, the purpose of punishment and prison sentences has shifted from rehabilitation, where prisons serve as a way to prepare individuals for return to society, to deterrence and incapacitation.¹⁸⁵ Evidence-based sentencing puts incapacitation at the front and center of the debate. Because recidivism is the focus of evidence-based sentencing, the idea is that the public will be protected while the defendant is imprisoned and thus incapacitated.

Actuarial risk assessments give a simple numerical value—a percentage with questionable accuracy—to answer what appears to the lay person to be the probability that a person will recidivate.¹⁸⁶ This does not provide information about when or how they will recidivate. This snapshot is provided to us after a plea or a trial, but before sentencing, and ignores the many factors that can impact recidivism after and during an individ-

¹⁸¹ Brief for the American Medical Association et al., as *Amici Curiae* Supporting Respondents at 5, *Roper v. Simmons*, 543 U.S. 551 (2005) (No. 03-633) (quoting L.P. Spear, *The Adolescent Brain and Age-Related Behavioral Manifestations*, 24 NEUROSCIENCE & BEHAV. REVS. 417, 421 (2000)).

¹⁸² See *Miller*, 567 U.S. at 483 (noting that mandatory punishment “disregards the possibility of rehabilitation even when the circumstances [(juvenile brain development and culpability)] most suggest it”).

¹⁸³ See Starr, *supra* note 1, at 816.

¹⁸⁴ See Oleson, *supra* note 8, at 1330–32.

¹⁸⁵ See Michelle S. Phelps, *Rehabilitation in the Punitive Era: The Gap Between Rhetoric and Reality in U.S. Prison Programs*, 45 L. & SOC’Y REV. 33, 34 (2011).

¹⁸⁶ Indeed, what qualifies as recidivism is hard to define. See *supra* subpart III.B (Prong I: Whether Evidence-Based Sentencing Has been Tested?).

ual's incarceration (including family ties¹⁸⁷ and prison education programs¹⁸⁸). Using a single percentage generated by a risk assessment tool does not reduce recidivism; rather, it provides a false feeling of creating a "just" criminal justice system, while failing to rehabilitate and help offenders transition back into society.

CONCLUSION

Despite its purported positive impact on the criminal justice system, evidence-based sentencing risks fooling judges and juries and further contributing to the overrepresentation of men of color and poor people in prisons. The problems with the creation of these models, namely a lack of replication, potentially unconstitutional use of certain factors, high false positive rates, and issues with G2i abstraction, should all create serious concerns for actors in and around the criminal justice system.

The *Daubert* test offers an analytical framework through which the validity and fit of evidence-based sentencing can be evaluated. As a general matter, evidence-based sentencing can and has been tested, but rigor should be applied to sample selection and methods. Peer review should be decided on a case-by-case basis and free from conflicts of interest. The high error rates of specific evidence-based sentencing tools are concerning, particularly with the inability to pinpoint what a large effect size is. Finally, the method of using logistic regression for individual prediction is a generally accepted statistical technique, but the weight to give such information has not been established. The Pennsylvania Commission on Sentencing's current proposal does not pass the first three prongs of *Daubert*, and general acceptance could be argued either way depending on the scope of analysis. Even if Pennsylvania's evidence-based sentencing regime did pass the *Daubert* test, the prejudicial nature of using sex and age would disallow evidence-based sentencing. Whether this sentencing technique would be appropriate without such factors remains an open question.

It will be hard to measure the impact of Pennsylvania's risk assessment tools when they go into effect on July 1, 2018. The two risk assessment tools for a future offense of "Any Crime" or

¹⁸⁷ Mark T. Berg & Beth M. Huebner, *Reentry and the Ties that Bind: An Examination of Social Ties, Employment, and Recidivism*, 28 JUST. Q., 382, 383 (2011).

¹⁸⁸ From class discussion through the Cornell Prison Education Program.

“A Crime Against a Person” may sometimes result in an offender being in different risk groups depending on the test applied.¹⁸⁹ The guidelines offer no explanation on how to reconcile these two outcomes or whether further information is needed if just one of the risk assessments results in a low or high risk. Measuring the outcome of seeking further information for low- or high-risk assessments may be difficult as the judge may order either “a PSI report that contains risk-needs-responsivity information” or “the preparation of an RNA or RNR assessment.”¹⁹⁰ After considering its methodological and constitutional shortcomings, it is unclear whether Pennsylvania’s new risk assessment tools will truly result in “smarter sentencing,” and whether we will know what it did at all.

¹⁸⁹ PROPOSALS PUBLISHED, *supra* note 85, § 305.9, Risk Assessment Summary (the Commission failed to label § 305.9 in the original report).

¹⁹⁰ *Id.* § 305.5(c).

